

# Imperial College London

IMPERIAL COLLEGE LONDON

DEPARTMENT OF LIFE SCIENCE

---

## Genomic Architecture of Parallel Ecological Divergence in *Littorina saxatilis*

---

*Author:*

Rui Zhang

*Supervisor:*

Matteo Fumagalli

m.fumagalli@imperial.ac.uk

*CID:*

01907894

*Co-supervisor:*

Francesca Raffini

f.raffini@sheffield.ac.uk

*Email:*

rui.zhang20@imperial.ac.uk

Roger K. Butlin

r.k.butlin@sheffield.ac.uk

*Date:* 8/26/2021

A THESIS SUBMITTED FOR THE PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE  
DEGREE OF MASTER OF SCIENCE AT IMPERIAL COLLEGE LONDON

SUBMITTED FOR THE MSc IN COMPUTATIONAL METHODS IN ECOLOGY AND EVOLUTION

## Declaration

I declare that this thesis is solely composed by my own work. The data was provided by co-supervisor Francesca Raffini and Roger K. Butlin from The University of Sheffield. The preliminary data processing was conducted by co-supervisor Francesca Raffini. I only executed data analysis and results visualization. Supervisor Matteo Fumagalli provided a workshop on how to use ANGSD v.0.934. Before August 2021, this thesis has not been published in any journal.

Signature:

A handwritten signature in brown ink that reads "Rui Zhang". The signature is written in a cursive style with a clear distinction between the two characters.

Date: 8/26/2021

# Contents

<b>Abstract</b>	<b>3</b>
<b>Keywords</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
<b>2 Method</b>	<b>7</b>
2.1 Sampling . . . . .	7
2.2 DNA sequencing and processing . . . . .	7
2.3 PCA . . . . .	7
2.4 Admixture analysis . . . . .	8
2.5 Folded one dimensional and two dimensional site frequency spectrum . . . . .	8
2.6 Population genetic differentiation and nucleotide diversity . . . . .	8
2.7 Linkage disequilibrium heat map and decay analysis . . . . .	9
<b>3 Results</b>	<b>10</b>
3.1 PCA . . . . .	10
3.2 Admixture analysis . . . . .	10
3.3 1d SFS . . . . .	13
3.4 2d SFS . . . . .	14
3.5 Summary statistics . . . . .	15
3.6 Linkage disequilibrium . . . . .	18
<b>4 Discussion</b>	<b>20</b>
4.1 PCA and admixture analysis . . . . .	20
4.2 Summary statistics of differentiation . . . . .	20
4.2.1 SFS . . . . .	20
4.2.2 Population genetic differentiation . . . . .	21
4.2.3 Limitation of FST . . . . .	21
4.3 Linkage disequilibrium . . . . .	22
4.4 Other possible mechanisms of parallel adaptation . . . . .	23
4.5 Future direction . . . . .	23
<b>5 Conclusion</b>	<b>24</b>
<b>Code availability</b>	<b>24</b>
<b>Acknowledgement</b>	<b>25</b>
<b>Appendix</b>	<b>30</b>

## <sup>1</sup> Abstract

<sup>2</sup> Parallel adaptation of *Littorina saxatilis* in heterogeneous environments has been researched by  
<sup>3</sup> many scientists all around the world, since it is a good subject to explore sympatric speciation  
<sup>4</sup> and natural selection. This species has evolved two ecotypes adapted to micro-habitats: crab  
<sup>5</sup> ecotype and wave ecotype. The aim of this research is to find underlying mechanisms of parallel  
<sup>6</sup> ecological divergence with gene flow and the evolvement of reproductive isolation with inter-  
<sup>7</sup> breeding. Low coverage whole genome sequencing and related data analysis tools were used on  
<sup>8</sup> 73 samples of snails and the research on hybrids is the emphasis.

<sup>9</sup>

<sup>10</sup> In the Spanish system of *Littorina saxatilis*, snails are divided into two distinct genetic clusters  
<sup>11</sup> instead of a cline in Sweden. The crab group and hybrid group are genetically closer and  
<sup>12</sup> they compose the first cluster, a cline. The wave group forms the second homogeneous cluster.  
<sup>13</sup> Also, there is an obvious directional introgression from wave into crab at a low rate which  
<sup>14</sup> destroys high Linkage Disequilibrium (LD) in crab ecotype. Differentiation analysis indicates  
<sup>15</sup> that most loci are polymorphic and not fixed suggesting that loci are under both divergent  
<sup>16</sup> selection and balancing selection. Usually, chromosome inversion candidates show higher FST  
<sup>17</sup> and LD compared with whole chromosome, while PBS and genetic diversity ( $\pi$ ) are more diverse.  
<sup>18</sup> The crab group has the highest LD and slowest LD decay rate which hints at stronger selection  
<sup>19</sup> and more polymorphic inversions. It is hard to find new chromosome inversions just from LD  
<sup>20</sup> heat map, but some chromosomes demonstrate earlier points dispersal in LD decay suggesting  
<sup>21</sup> short inversions. In conclusion, this research found some new chromosome inversions candidates  
<sup>22</sup> to be verified and tried to explain the process of ecological divergence after the phase of a cline  
<sup>23</sup> in terms of genetic structure and selection.

## <sup>24</sup> Keywords

<sup>25</sup> Whole genome sequencing, parallel evolution, local adaptive, genomic architecture, gene flow,  
<sup>26</sup> reproductive isolation

## 27 1 Introduction

28 How speciation occurs is a fundamental biological problem. There are three traditional geo-  
29 graphic modes of speciation: allopatric speciation, parapatric speciation, and sympatric specia-  
30 tion (Fitzpatrick et al. 2008). Sympatric speciation is fascinating due to its difficulty in theory  
31 explanation and definition. Parallel ecological divergence is one of the most interesting subjects  
32 of sympatric speciation. A monomorphic population could split into two coexisting phenotypic  
33 clusters (ecotypes) under directional selection of ecological environments explained by the theory  
34 of adaptive dynamics (Dieckmann & Doebeli 1999). It is well established that divergent natural  
35 selection in heterogeneous ecological environments is likely to be an impetus for local adaptation  
36 and subsequent reproductive isolation even speciation (Schluter 2009). Repeated occurrence of  
37 ecological divergence in discrete populations in contrasting habitats gives an opportunity to ex-  
38 plore speciation mechanisms. Moreover, how reproductive isolation evolves with interbreeding,  
39 which means the presence of gene flow impeding speciation, is more informative to understand  
40 initial steps of speciation than the study on already separated species.

41

42 *Littorina saxatilis* is a common rocky-shore gastropod with short life and high lifetime dispersal  
43 (Reid 1996). This species has a strong adaptation ability. Even small patches of local habitat  
44 could promote a distinct ecotype under biotic or abiotic pressure (Johannesson 2016). It has  
45 evolved two ecotypes under the selection of wave exposure and crab predation in many loca-  
46 tions such as northwest Spain, west coast of Sweden, and northeast coast of England (Butlin  
47 et al. 2008). Wave ecotype is smaller and has a thinner shell with a larger aperture, while crab  
48 ecotype is bigger and has a thicker shell with a smaller aperture (Johannesson et al. 2010).  
49 The morphology and behavior differences between two ecotypes are inheritable although with  
50 a minor phenotypic plasticity (Galindo et al. 2009). *Littorina saxatilis* has developed partially  
51 reproductive isolation in the face of hybridization (Grahame et al. 2006). This proposition is  
52 supported by that gene flow in contact zones was 10-30 % of gene flow within ecotypes and that  
53 the genetic relationship between crab ecotype and wave ecotype in the same location is closer  
54 than the genetic relationship between the same ecotype in different locations (Panova et al.  
55 2006). Also, Beaumont (Beaumont 2010) got strong support for parallel and local divergence  
56 in this species instead of old allopatric divergence of ecotypes followed by secondary overlap  
57 and gene flow using approximate Bayesian computation (ABC) approach. Therefore, it is a  
58 preferable model system to study the underlying mechanisms, particularly genomic architecture  
59 variation, of parallel adaptive divergence with the existence of gene flow, which is sympatric or  
60 parapatric speciation instead of allopatric speciation.

61

62 There are many studies focusing on the history of *Littorina saxatilis* dispersal, colonization,  
63 formation of ecotypes, and evolution of reproductive barriers between two ecotypes under gene  
64 flow. This species is inferred to survive from the last glacial period in a northern latitudes  
65 refugia based on phylogeographic data (Panova et al. 2011). Later colonization seemed to be  
66 achieved by rafting of single females whose brood pouch carried hundreds of embryos sired by  
67 even over 20 males. This feature allows a single female to be a founder group with diversified

68 genetic variation and a new population is established rapidly once releasing embryos (Rafajlovic  
69 et al. 2013). Crab is regarded as the driving force of crab ecotype formation. After snails' colo-  
70 nization, the predator of snails arrived at the same location later due to their higher minimum  
71 temperature requirement compared with snails. In the ABC model established by Beaumont  
72 (Beaumont 2010), ecotypes separation is a relatively recent incident. The ecotype separation  
73 time is estimated at only around 10 % of local population age. Furthermore, the formation of  
74 ecotypes occurs instantaneously (<1000 generations) rather than gradually. Although ecotypes  
75 have formed, individuals in two ecotypes are still able to interbreed with each other even from  
76 different locations (Hollander et al. 2005). The survival rate of hybrids in the contact zone  
77 is higher than parental ecotypes or on the same level which means hybrid superiority exists  
78 (Rolan-Alvarez et al. 1997). It promotes the first step of speciation, ecotype formation because  
79 the primary contact zone is easy for individuals to expand the population from one ecotype over  
80 the contact zone to the other microhabitat. However, it impedes the completion of speciation  
81 due to continuous gene flow.

82

83 Gene flow between crab and wave ecotypes never disappear although in the contact zone gene  
84 flow is impeded and less than gene flow within ecotypes (Panova et al. 2006). How *Littorina*  
85 *saxatilis* maintain divergence in the face of gene flow which is regarded as genetic recombination  
86 counteracting divergence? Strong enough selection pressure could overcome this kind of homoge-  
87 nizing effects and genomic architectures might resist gene flow by impeding gene recombination  
88 (Smadja & Butlin 2011). Hybrid zone analysis has inferred patterns of selection in space of  
89 *Littorina saxatilis* that crabs represent a strong selection pressure during ecotypes separation  
90 (Westram et al. 2018). Chromosomal inversions are regarded as reservoirs and vehicles for rapid  
91 parallel divergence in *Littorina saxatilis* (Morales et al. 2019). In addition, assortative mating  
92 of ecotypes is obvious in both field and laboratory experiments. Individuals are more likely to  
93 mate with the same ecotype individuals by means of following trails and longer mating time  
94 (Hollander et al. 2005). In brief, the partial speciation of this species is considered to be evolved  
95 by divergent selection of heterogeneous environments, habitat choice, assortative mating, and  
96 genomic background (chromosome rearrangement).

97

98 In our research site, Spanish, ecotypes are distributed over vertical shore gradients. Hybrids  
99 are a minority compared with parental ecotypes in the intermediate habitat (Johannesson et al.  
100 1993) with isolation indexes of 0.5–0.9 in Spain tested from mate choice experiments (Johan-  
101 nesson et al. 1995). The Spanish system of *Littorina saxatilis* ecological divergence is found  
102 much older than Swedish and British systems based on mitochondrial DNA lineages (Panova  
103 et al. 2011). Furthermore, Morales (Morales et al. 2019) found that in transects spanning the  
104 crab-hybrid-wave axis, the comparisons between crab ecotype and wave ecotype showed more  
105 significant divergence in Spain than in Sweden. However, the comparisons only focus on crab  
106 ecotype and wave ecotype snails themselves and ignore hybrids. Hybrid of two ecotypes could  
107 offer more information than distinct population comparisons. For instance, cline analysis and hy-  
108 brid zone analysis in Sweden (Westram et al. 2018) were used to detect non-neutral SNP (single

109 nucleotide polymorphism) and find a genotype-phenotype-environment association. Therefore,  
110 this project will take hybrid whole genome sequencing into account.

111

112 Then what kind of sequencing method and subsequent data analysis should be used in this  
113 research? Low coverage whole genome sequencing (lc WGS) is regarded as a cost-efficient ap-  
114 proach to catch low frequency variation in many individual samples of population (Gilly et al.  
115 2018). The development of calling algorithms makes it possible to call SNP and estimate allele  
116 frequency accurately from low depth sequencing data. Analysis of next generation sequencing  
117 data (ANGSD) is a software designed for next generation sequencing data with data analysis  
118 methods of taking genotype uncertainty into account instead of calling genotypes directly (Ko-  
119 rneliussen et al. 2014). It is especially suitable and useful for low and medium depth data.  
120 Therefore, it is a good sequencing data analysis tool for lc WGS data of population genomics.  
121 ANGSD has been used to explore underlying genetic mechanisms of parallel phenotypic evolu-  
122 tion, local adaptation, history of evolution, and so on (Wilder et al. 2020) (Therkildsen et al.  
123 2019) based on population genetic sequencing data. Furthermore, linkage disequilibrium is a  
124 good indicator to explore population genetic history including selection, domestication, chromo-  
125 some rearrangement, and so on. Also, LD decay analysis could reveal population recombination  
126 and selection history (Zhang et al. 2018). All of these measures help us to explore population  
127 structure, natural selection signature, and history of genetic changes.

128

129 In this research, I expect to find distinct genetic groups from principal components analysis  
130 (PCA), special population genetic differentiation, and nucleotide density pattern in chromo-  
131 somes, detect new inversions candidates and infer parallel adaptive evolution history of *Littorina*  
132 *saxatilis* in Spain. Combining all above, I try to explain the underlying mechanism of parallel  
133 ecological divergence and the process of speciation facing gene flow.

<sup>134</sup> **2 Method**

<sup>135</sup> Samples collection, DNA sequencing and processing were conducted by stuff in the University  
<sup>136</sup> of Sheffield. Also details of these process were provided by co-supervisor Francesca and Roger.

<sup>137</sup> **2.1 Sampling**

<sup>138</sup> *Littorina saxatilis* snails were sampled in March 2017 on the west coast of Galicia in Spain:  
<sup>139</sup> Centinela (hereafter ER\_EA1; N 42° 4' 38.06", W 8° 53' 47.47"). The sample consisted of  
<sup>140</sup> approximately 600 snails from a 5 m wide band and separated by 0-5 m, stretching from the  
<sup>141</sup> upper limit of the species distribution in the splash zone to its lower limit close to low water of  
<sup>142</sup> spring tides. Sampling was denser in the lower part of the shore where hybridization between  
<sup>143</sup> 2 ecotypes described before was expected (Galindo et al. 2013). Snails were sampled by hand,  
<sup>144</sup> haphazardly, as far as possible without reference to phenotype but aiming to avoid juveniles.  
<sup>145</sup> All of snails were stored in individual tube, moistened with seawater and kept at 4 °C until  
<sup>146</sup> phenotyping was complete, then the head and foot of each snail was preserved in 100% ethanol.

<sup>147</sup> **2.2 DNA sequencing and processing**

<sup>148</sup> Low coverage whole genome sequencing (lcWGS) was conducted for 73 snails from ER\_EA1,  
<sup>149</sup> chosen to cover the full sampling range. DNA was extracted using the modified CTAB protocol  
<sup>150</sup> described by Panova (Panova et al. 2016). Library preparation (in-house high-throughput gDNA  
<sup>151</sup> library prep) and sequencing (Illumina HiSeq4000, 150bp PE) were carried out by The Oxford  
<sup>152</sup> Genomics Centre with target coverage 3x based on the estimated genome size of 1.35Gb. Samples  
<sup>153</sup> were sequenced in eight lanes, for a total of 16 lanes for each individual. Raw reads were trimmed  
<sup>154</sup> with Trimmomatic v.0.38 (<https://github.com/usadellab/Trimmomatic>) to remove the Illumina  
<sup>155</sup> adapters, mapped to the *L. saxatilis* draft reference genome (Swedend, crab; 1.35Gb) (Westram  
<sup>156</sup> et al. 2018) using bwa mem v.0.7.17 (<https://github.com/lh3/bwa>) and default settings, and fil-  
<sup>157</sup> tered by base and map quality 20 with Samtools v.1.7 (<https://github.com/samtools/samtools>).  
<sup>158</sup> PCR duplicates and PE overlap were removed with Picard v.2.0.1 ([https://github.com/broadin-  
<sup>159</sup> stitute/picard](https://github.com/broadinstitute/picard)) and bamUtil v.1.0.15 (<http://genome.sph.umich.edu/wiki/BamUtil>), respectively.  
<sup>160</sup> Bam files belonging to the same individual were sorted and merged with Samtools.

<sup>161</sup> **2.3 PCA**

<sup>162</sup> After getting bam files, I used software ANGSD v.0.934 (Korneliussen et al. 2014) (<http://www.popgen.dk/angsd/index.php/ANGSD>) to estimate imputed genotype probabilities from mapped  
<sup>163</sup> reads. ANGSD is a software designed for analyzing next generation sequencing data, especially  
<sup>164</sup> low and medium depth data due to its feature of taking genotype uncertainty into account. In  
<sup>165</sup> order to explore the population structure of *Littorina saxatilis* in Spain, I performed principal  
<sup>166</sup> components analysis (PCA) and admixture analysis. At first, I used ANGSD to perform SNP  
<sup>167</sup> calling, estimate genotype likelihood in beagle format and get a list of SNP position as input  
<sup>168</sup> files to estimate covariance matrix. Both ANGSD and PCAngsd could infer covariance matrix  
<sup>169</sup> of *Littorina saxatilis* 73 samples. Here I chose PCAngsd v.1.02 (Meisner & Albrechtsen 2018)

171 (<http://www.popgen.dk/software/index.php/> PCAngsd) and then used R v.4.0.3 to perform  
172 PCA using eigen function and extract eigenvectors to plot PC1 vs PC2. I also used individuals  
173 distance along crab-hybrid-wave axis data as color information to show population structure  
174 along crab-hybrid-wave axis.

## 175 **2.4 Admixture analysis**

176 Another method to interpret population structure is admixture analysis which showed genome-  
177 wide admixture proportions for every individual. ngsAdmix (<http://www.popgen.dk/software/in->  
178 dex.php/NgsAdmix) (Skotte et al. 2013) used beagle format genotype likelihood file to infer ad-  
179 mixture proportions of ancestry clusters for every individual. Then I used R to draw admixture  
180 proportions bar plot along crab-hybrid-wave axis. After getting admixture analysis result, I  
181 could classify 73 individuals into 3 groups: crab, hybrid and wave according to individual an-  
182 cestry clusters contribution. Given ecotype information, I could draw new PCA plot in 3 colors  
183 showing population structures of every ecotype population.

## 184 **2.5 Folded one dimensional and two dimensional site frequency spectrum**

185 In order to detect selection signatures on genome and understand patterns of divergence and  
186 differentiation across genome, I calculated several summary statistics from low-depth NGS data.  
187 The first step was estimating sample allele frequencies (SAF) posterior probabilities by ANGSD  
188 for each population separately. Then I used program realSFS of ANGSD to calculate Site  
189 Frequency Spectrum (SFS) which records the proportions of sites at different allele frequencies.  
190 Here I calculated folded SFS without outgroup species defining ancestral state. After getting  
191 one dimensional folded SFS, I drew bar plot for every population removing the first value of 1D  
192 folded SFS which represented the expected number of sites with derived allele frequency equal to  
193 0 due to its relatively too big value to guarantee other values clear interpretation. Next step was  
194 estimating joint SFS between 2 populations (2D folded SFS) which is useful to infer divergence  
195 process of populations and estimate FST and PBS as prior information. As before calculating  
196 1D folded SFS, I used program realSFS to infer 2D folded SFS between every pair of 3 ecotype  
197 populations. Similarly, I could use R to illustrate these flatten matrixes in heat map format.

## 198 **2.6 Population genetic differentiation and nucleotide diversity**

199 Here I chose FST and population branch statistic (PBS) as indicator of allele frequency differen-  
200 tiation. ANGSD helps us calculate FST and PBS from sample allele frequencies likelihoods (.saf  
201 files), avoiding genotype calling. The first step was computing per-site FST indexes by realSFS  
202 using SAF files and 2D folded SFS as input. The output file of this stage is (a) and (a+b) values  
203 for three FST comparisons for every SNP sites. I could use R to calculate FST value for every  
204 SNP and every pairwise of 3 populations and then drew FST in every chromosome and every  
205 pairwise populations in the order of position. The next step was performing a sliding-window  
206 analysis with setting window size as 10kb and step size as 1kb by realSFS. Similarly I drew  
207 window PBS value in every chromosome and population.

208

209 Furthermore, in order to know whether allele frequency differentiation increase is related to  
210 the change of nucleotide diversity, I used realSFS saf2theta function and program thetaStat of  
211 ANGSD to calculate thetas for each site and in windows taking SAF and SFS as input files and  
212 prior information. Pairwise theta estimation ( $tP$ ) divided by numbers of sites in the window  
213 ( $nSites$ ) can be used to estimate window-based pairwise nucleotide diversity ( $\pi$ ).

214 **2.7 Linkage disequilibrium heat map and decay analysis**

215 Linkage disequilibrium is a useful indicator to illustrate correlation of pairwise locis. Here  
216 are various different measures of LD such as pairwise  $r^2$ , pearson estimation of  $r^2$ , D from EM  
217 algorithm and so on. Here I choose pairwise  $r^2$  (coefficient of correlation) because it is very useful  
218 when analyzing biallelic markers such as SNPs and independent of sample size (Devlin and Risch  
219 1995). I subset beagle files and maf files of every chromosome at first, and then sample these  
220 SNP sites to guarantee total SNP number of every chromosome is around 10000 as instruction  
221 of ngsLD software. ngsLD v.1.1.1 (Fox et al. 2019)(<https://github.com/fgvieira/ngsLD>) is a  
222 software designed for NGS data taking the uncertainty of genotype's assignation into account.  
223 After using this software calculating LD among SNPs in every chromosome, I used R to draw  
224 LD heat map in whole chromosome and partial chromosome inversion candidates. Then I used  
225 a modified R script offered by ngsLD to analyze LD decay and draw LD decay line for every  
226 chromosome and population.

227 **3 Results**

228 **3.1 PCA**

229 The covariance matrix of principal components analysis (PCA) could be estimated by single-  
230 read sampling in ANGSD. With proper parameters settings, there are 12,817,872 variant sites  
231 (SNP) called from 73 individuals bam files. Here I used PCAngsd to take genotype likelihoods  
232 in beagle format as input and then infer covariance matrix which is used in Figure 1.

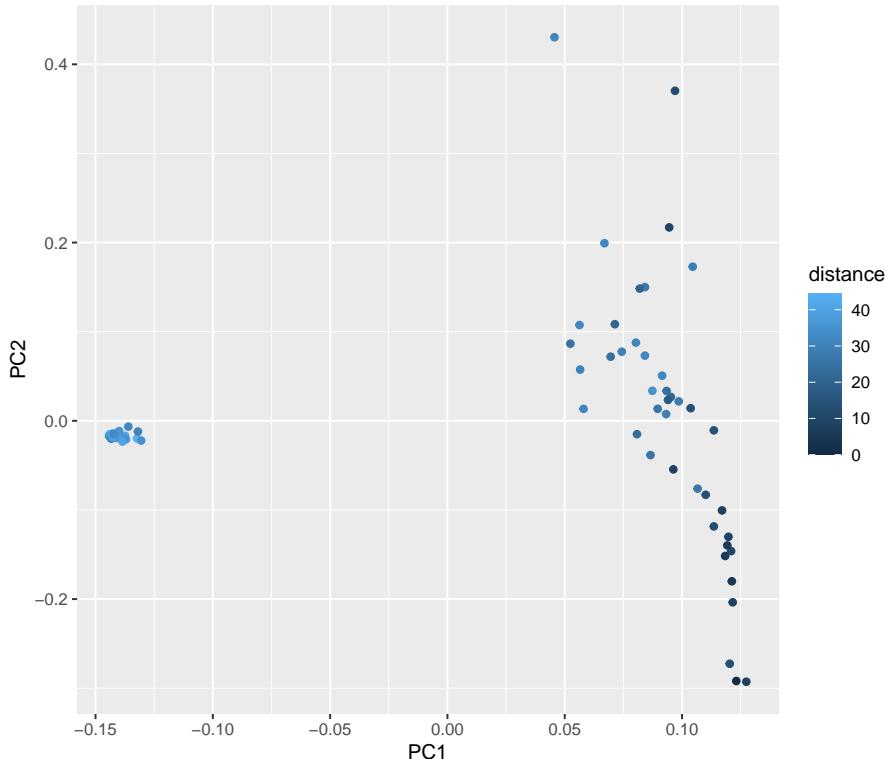


Figure 1: PCA of 73 *Littorina saxatilis* individuals

Dark blue presents crab ecotype and light blue presents wave ecotype. The more dark means the sampling individual is closer to crab ecotype along the crab-hybrid-wave axis. V2 is the first principal component and V3 is the second principal component.

233 As Figure 1 showed, there are many individuals points squeezing together in a little circle region  
234 and the color of these individuals are light blue which means these individuals are likely to  
235 be wave ecotype snails. On the other side of Figure 1, points are darker and comparatively  
236 scattering. Furthermore, deep dark blue points are usually lower than medium blue points. So  
237 far, I could only know that there are 2 genetic distinct clusters along crab-hybrid-wave axis.

238 **3.2 Admixture analysis**

239 In order to infer every individual's genome-wide admixture proportions, I used ngsAdmix to  
240 estimate admixture proportions from genotype likelihoods and ancestry clusters number setting.

241

242 As Figure 2A showed, right individuals are mostly wave ecotypes while left individuals are mostly

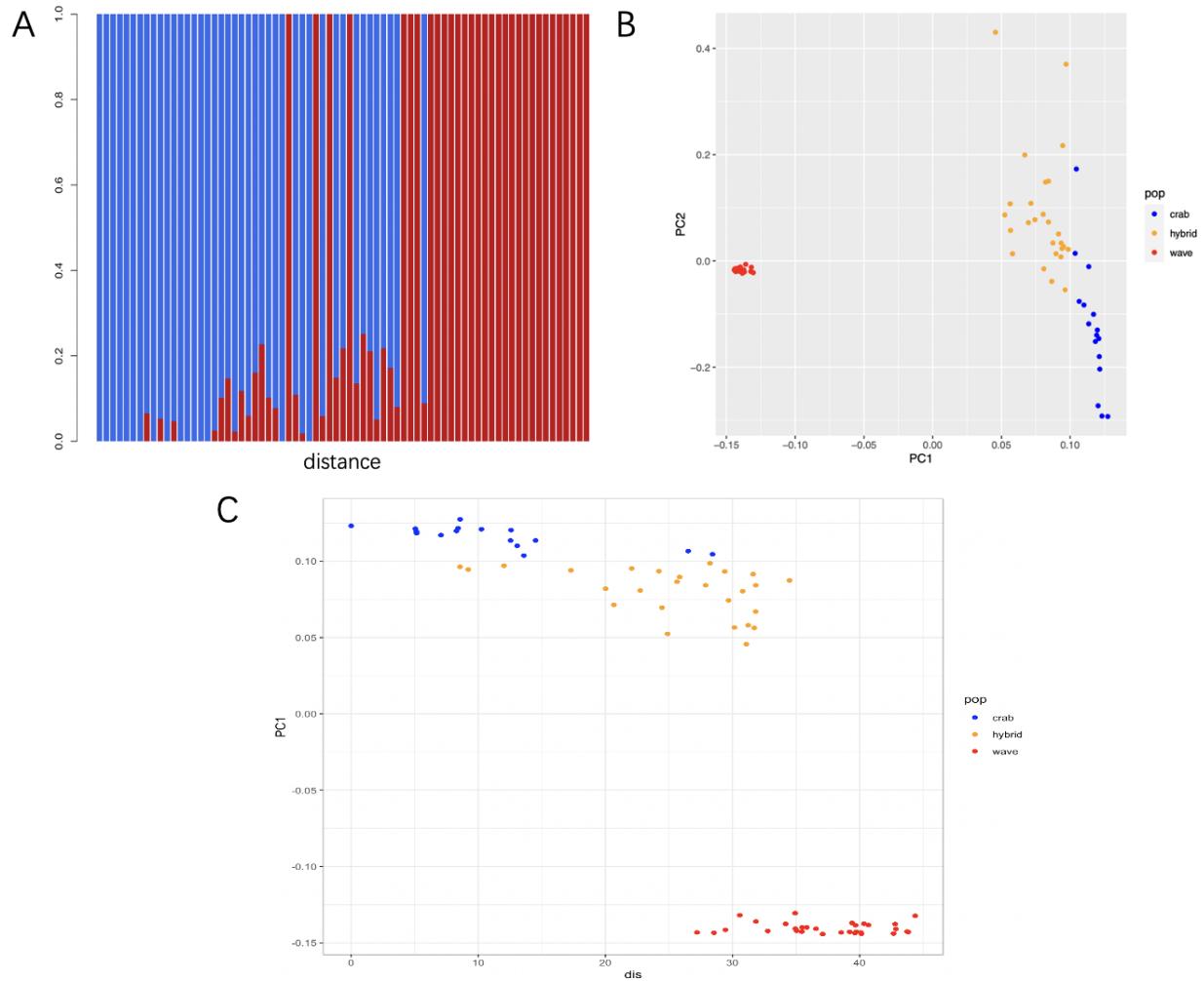


Figure 2: Admixture analysis of 2 ancestry clusters and new PCA based on admixture analysis  
**A.** The figure showed individuals admixture analysis results along crab-hybrid-wave axis from crab to wave ecotype. Blue represents crab ecotype and red represents wave ecotype.

**B.** Based on admixture analysis, individuals could be divided into 3 groups: wave, hybrid and crab. Given the group information, PC1 vs PC2 of principal components analysis is shown in 3 colors.

**C.** Based on admixture analysis, 3 groups named wave, hybrid and crab could be plotted in a PCA between distance of crab-hybrid-wave axis and PC1.

243 crab ecotypes which corresponds with distance data along crab-hybrid-wave axis. However, in  
244 the middle part of Figure 2A, there are crab ecotype, wave ecotype and hybrid individuals,  
245 which means in the contact zone, 2 ecotypes and their hybrid offspring are mixed. Furthermore,  
246 hybrid individuals can also enter pure wave and pure crab ecotype zone. From the result of  
247 admixture analysis, I could divide individuals into 3 groups: crab ecotype, wave ecotype and  
248 hybrid with 16, 31 and 26 individuals separately. When I used group information, I could redraw  
249 PCA plot shown as Figure 2B. All of wave ecotype individual points squeezed together just as  
250 Figure 1. Crab ecotype individuals and hybrid individuals are all on the other side, far away  
251 from wave ecotype, which means crab group and hybrid group are much more close compared  
252 with wave group. Meanwhile, hybrid group are higher in the second principal component than  
253 crab group, although it is difficult to distinguish hybrid group and crab group only using the  
254 first two principal components.

255

256 In Figure 2C, all wave individuals are around -0.14 in PC1, all crab individuals are higher  
257 than 0.10 while hybrid individuals are lower than 0.10 in PC1. However, only using the line of  
258 PC1=0.10 is not sufficient to split the higher cluster into crab and hybrid two isolated popula-  
259 tions. Therefore, I would like to regard that there are 2 distinct groups, one is a cline within  
260 crab and hybrid group, the other is wave group.

261 3.3 1d SFS

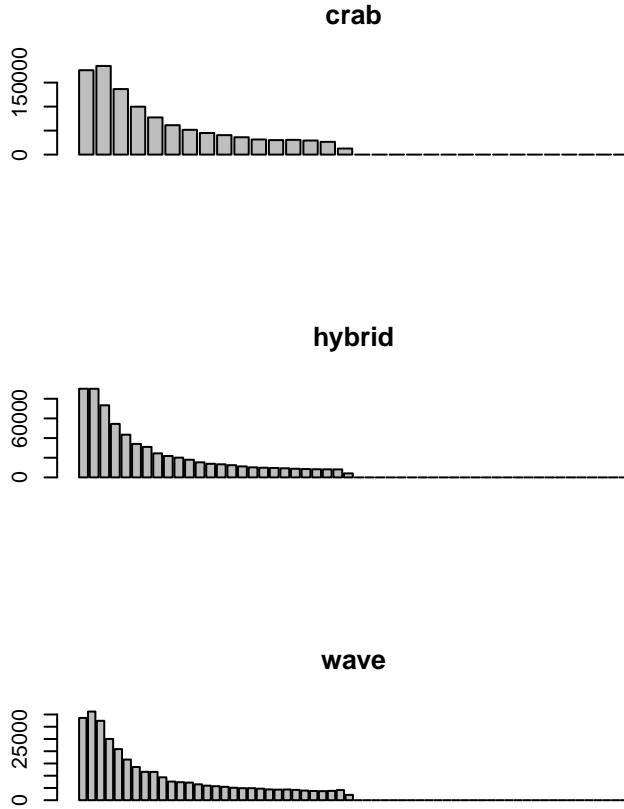


Figure 3: One-dimensional folded Site Frequency Spectrum (SFS)

One-dimensional folded SFS of each group is shown in 3 bar plots from allele frequency 1 to the number of individuals in the group

262 Site frequency spectrum records the number of sites at different allele frequencies. One-dimensional  
263 folded SFS only contains allele frequencies in one group. In Figure 3, sites will be recorded in  
264 SFS only if this site is caught in every individual of the group, also I have removed all of the  
265 sites without alleles due to its comparatively huge value. As the bar plots showed, there are  
266 less sites of allele frequency 1 than allele frequency 2 in crab and wave group. It is unexpected  
267 but it might be caused by distortion effect on SFS of polymorphic inversions and SNP filter  
268 condition. The sites number of crab group is bigger than hybrid and wave group. It could be  
269 explained by smaller population size of crab group which means it is easier for crab group to  
270 detect the same site in every individual. Furthermore, after calculating the proportion of SNP  
271 sites in every group, crab group SNP proportion is 0.0193, wave group SNP proportion is 0.0185,  
272 hybrid group SNP proportion is 0.02266. Hybrid group has the highest SNP proportion maybe  
273 because hybrid group contains SNP both in crab group and wave group.

274 3.4 2d SFS

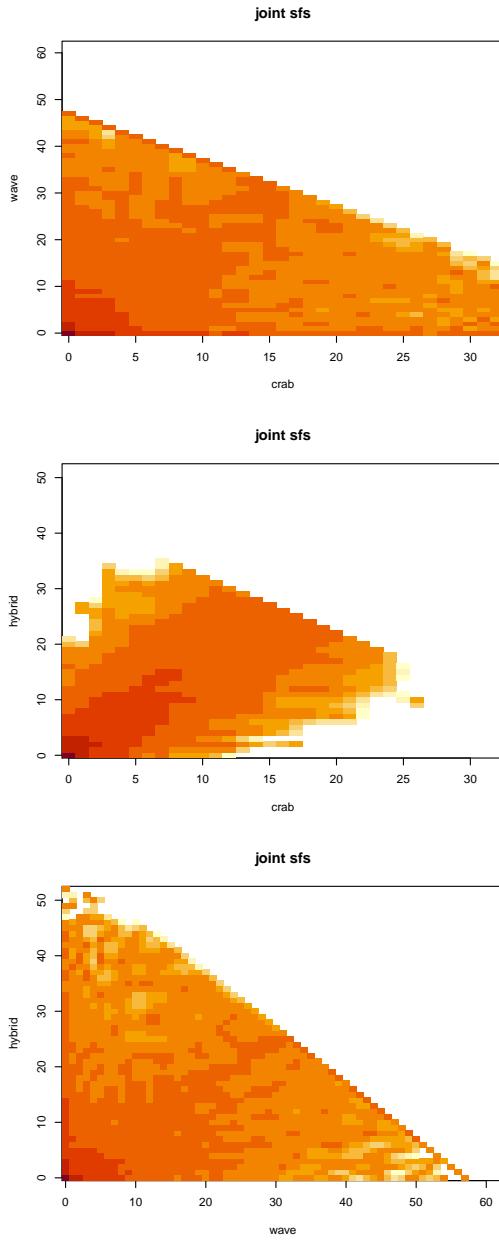


Figure 4: Two-dimensional folded Site Frequency Spectrum (SFS)

Two-dimensional folded SFS of each pair of groups are shown in 3 heat plots from allele frequency 0. Darker color means more sites and lighter color means less sites.

275 Joint SFS between 2 populations could be used to infer divergence process. From Figure 4, it  
 276 is clear that crab and hybrid group are closer because there are more darker color squares on  
 277 the diagonal compared with the other 2 heat plots. Usually colored squares of folded 2D-SFS  
 278 are all under one of the diagonals. However in Figure 4, it is not the case. This situation might  
 279 be resulted by different population sizes. 2D-SFS could also be used as prior information for  
 280 estimating FST and PBS (population branch statistic) which will be discussed later.

281 3.5 Summary statistics

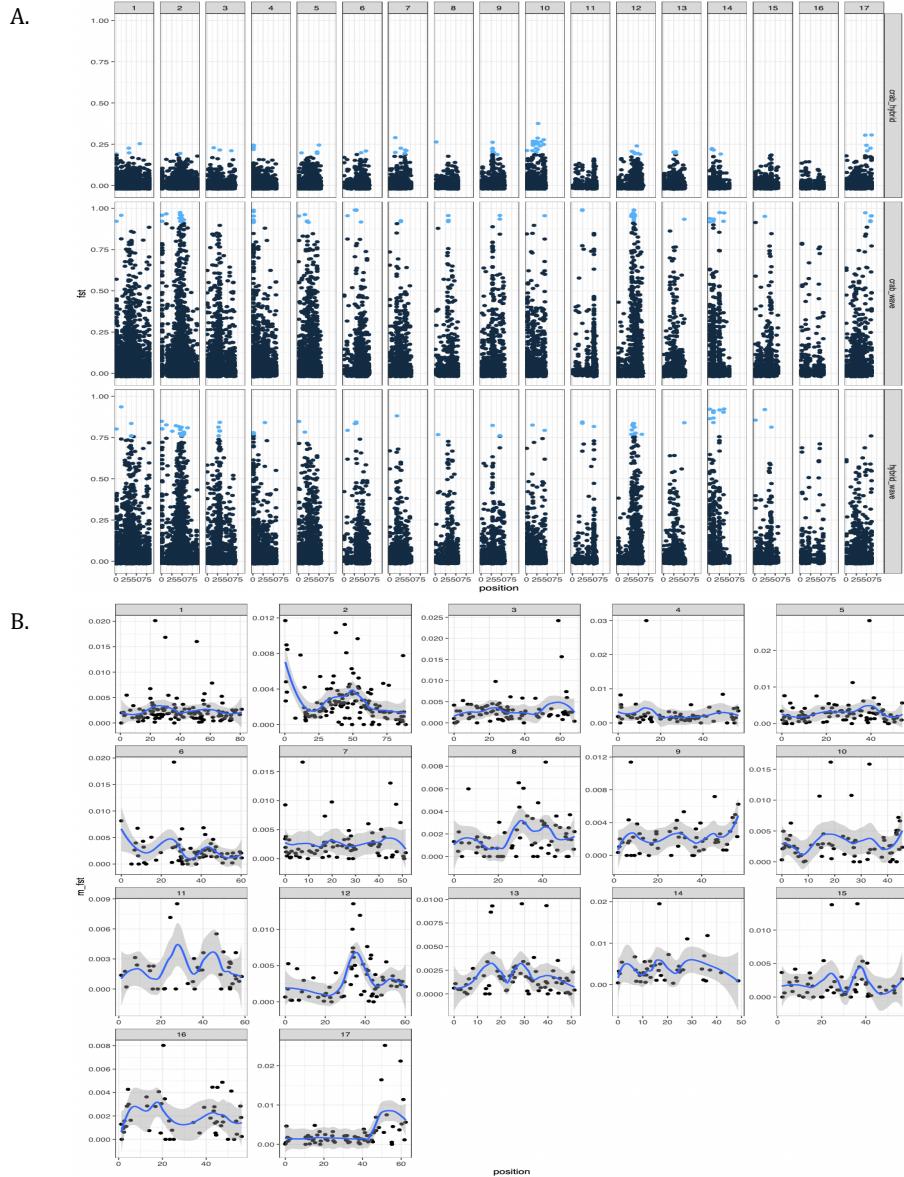


Figure 5: FST value among 17 chromosome and three pairwise populations

A. FST value scatter plot of every chromosome and pairwise populations with blue points representing outliers.

B. Mean FST value between crab and wave ecotype in every position on chromosome along with smooth lines showing trend.

282 There are some big blocks showing high FST value, high PBS value and low pi value, which  
 283 are all signatures of positive selection, especially in LG 1, 2, 4, 6, 9, 12, 14, 17. From Figure  
 284 5A, I can find that population genetic differentiation of crab and wave is the largest. Also, it is  
 285 clear that hybrid and crab are closer due to their low FST value in all chromosomes. Figure 5B  
 286 exhibits mean FST value of crab and wave population showing general high FST value in the  
 287 specific position such as crest part of LG2, 6, 12, 17. Population branch statistics of hybrid is  
 288 the lowest which is not difficult to understand because hybrid contain genetic components from

289 crab and wave ecotype. PBS of wave is the highest which is another proof of closer relationship  
290 between crab and hybrid. Pairwise nucleotide in Figure 2 in appendix can not offer too many  
291 effective information. The pi value looks nearly all the same in every chromosome and every  
292 ecotype.

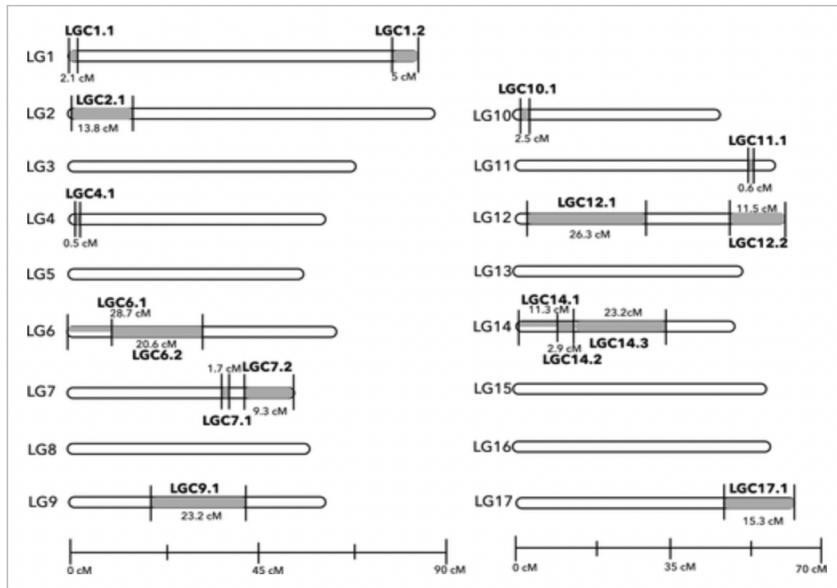
293

294 Faria et al. (Faria et al. 2019) found some chromosome inversion candidates as Figure 6A  
295 shown. In these regions, FST and PBS are usually relatively higher than other regions, while pi  
296 are usually relatively lower than other regions. In order to interpret differences between these  
297 chromosomal rearrangement candidates and other regions, I can also calculate mean value of  
298 FST, PBS, pi and linkage disequilibrium (LD, which will be discussed later) in these regions  
299 and whole chromosome. Table S1 in supplementary material of appendix gives exact number of  
300 these statistics values which is more helpful to some extent. Also, a scatter plot Figure 6B has  
301 been drawn to interpret differences between candidates and whole chromosome more intuitively.  
302 From Figure 6B, I can come to conclusion easily that generally FST and LD value of chromosome  
303 inversion candidate regions are higher than whole chromosome. However, PBS and pi value of  
304 chromosome inversion candidates region is sometimes lower and sometimes higher than whole  
305 chromosome. Statistics of LG1.2, 2.1, 4.1, 6.1, 6.2, 9.1, 14.3 meet the expected result of higher  
306 PBS and lower pi in inversion candidates regions. LG1.1, 7.1, 7.2, 17.1 show higher PBS and  
307 higher pi while LG 10.1, 11.1, 12.1, 12.2, 14.1, 14.2 show lower PBS and higher pi.

308

309 However, there are also some chromosome regions showing high FST and PBS and low pi, but  
310 are not chromosomal rearrangement candidates, such as middle part of LG 2, 8, 11 and 15.  
311 They will be discussed later.

A.



B.

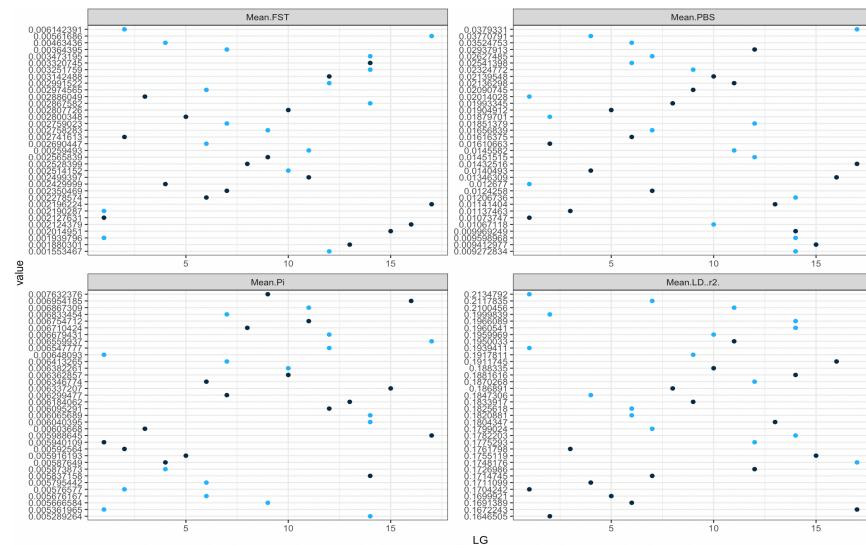


Figure 6: Mean statistics of inversion candidates and whole chromosome

A. Chromosomal rearrangement candidates inferred by Faria et al.(Faria et al. 2019)

B. Blue represents inversion candidates mean statistics while black represents whole chromosome mean statistics. The statistics here include mean FST of crab and wave ecotype, mean PBS of crab ecotype, mean pi of crab ecotype and mean LD of crab ecotype.

### 312 3.6 Linkage disequilibrium

313 Linkage disequilibrium estimates nonrandom association of alleles at different loci. Natural  
 314 selection, mutation, non-random mating, gene flow and population size can all affect LD. Here  
 315 I choose pairwise r<sup>2</sup> (coefficient of correlation) because it is very useful when analyzing biallelic  
 316 markers such as SNPs and independent of sample size (Devlin and Risch 1995). All of the LD  
 317 heat map and LD decay analysis regression plots are included in supplementary material of  
 318 appendix. Here Figure 7 only gives an example on chromosome 14. In 17 chromosomes, crab  
 319 population always has the highest LD value. Hybrid usually has medium LD value and wave  
 320 has the lowest LD value.

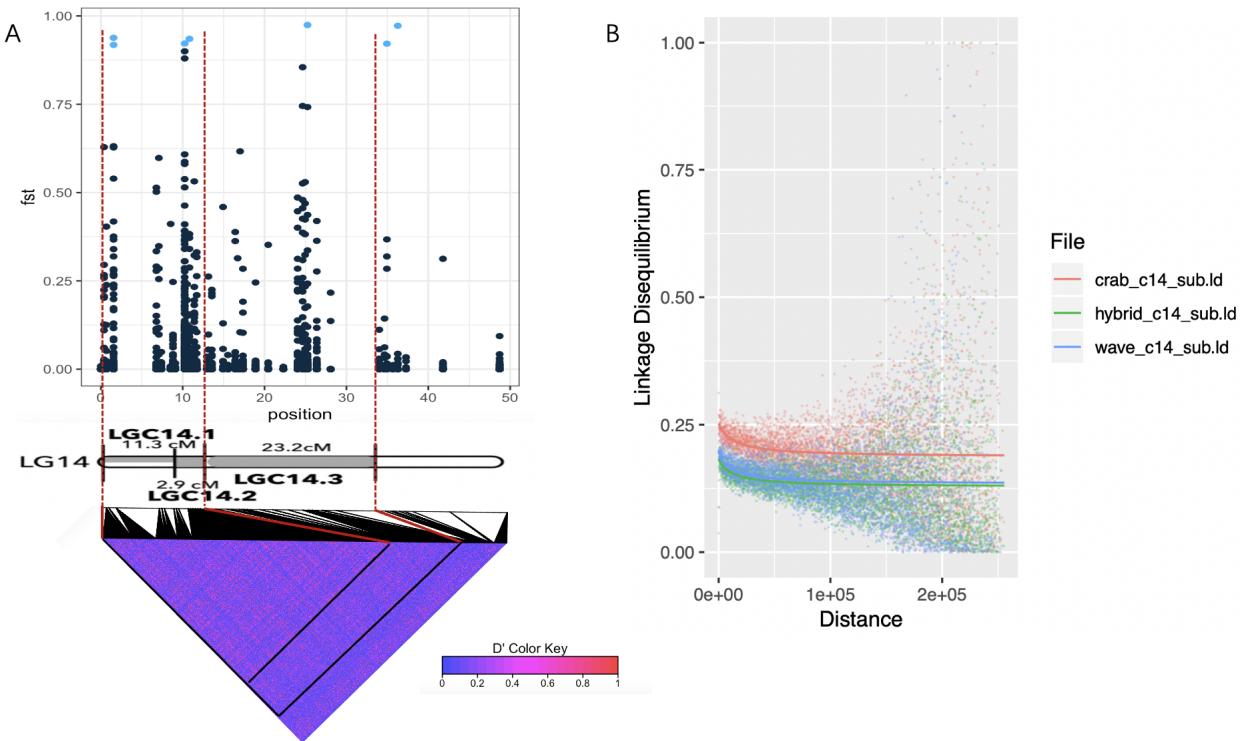


Figure 7: LD heat map and LD decay of chromosome 14

A. FST value with outliers mapping to chromosome inversion candidates and LD heatmap of crab population on chromosome 14.

B. LD decay analysis of chromosome 14

321 LD heat map among the whole chromosome is hard to get useful information because usually  
 322 all regions look the same. However in LG 14, the heat map between the first two red dot-  
 323 ted lines has a clear higher LD value block which maps to LG14.1 and LG14.2. The region  
 324 between the second and the third red dotted lines is also more pink than white region on chro-  
 325 mosome 14. Partial LD heat map in chromosome inversion candidates is more clearly showing  
 326 higher overall LD value compared with other regions. Exact mean LD values are shown in Table  
 327 S1 that mean LD values of chromosomal rearrangement candidates are higher than other region.

328

329 For LD decay analysis, every single point is mean LD value of a 50bp size window. The fitting  
 330 line decrease quickly approaching a horizontal line. As the distance increases, points become

<sup>331</sup> more dispersed and some chromosomes disperse more early such as LG1 and LG4. Crab usually  
<sup>332</sup> has the slowest LD decay rate in 17 chromosomes while hybrid ecotype population LD decay  
<sup>333</sup> rate is usually the fastest, which might be resulted by higher population genetic diversity gained  
<sup>334</sup> from crab and wave ecotype population.

335 **4 Discussion**

336 Here I try to explore mechanism of parallel adaptation of *Littorina saxatilis* facing gene flow.  
337 The genetic basis of parallel adaptation is supposed to depend on underlying genomic architecture  
338 (Yeaman 2013) such as the number, size and additivity of loci and nonrandom arrangements  
339 order in the genome. Chromosome inversion, a kind of chromosome rearrangement, has been  
340 researched widely. It is considered to be able to suppress recombination and serve as reservoirs  
341 for adaptive standing variation which might be the reason for fast parallel adaptation. Heteroge-  
342 neous environments induce directional selection on loci and create divergence on chromosomes.  
343 However, selection and chromosome rearrangement are not enough, hybridization counteracts  
344 divergence by the way of gene flow between 2 ecotypes. To explore the role of chromosome in-  
345 version, divergent selection of heterogeneous environments and hybridization during the initial  
346 stages of the speciation continuum in *Littorina saxatilis*, I used genetic analysis of sympatric  
347 populations of 2 ecotypes including intermediate individuals, hybrids.

348 **4.1 PCA and admixture analysis**

349 There are 2 clusters in PCA. The first group is more homogeneous and regarded as wave ecotype.  
350 The second group is more heterogeneous and regarded as a cline within crab and hybrid group.  
351 Furthermore, from the admixture analysis, I could find obvious directional introgression from  
352 wave into crab at a low rate.

353  
354 In Sweden admixture analysis conducted by Sheffield University stuff, there were some individ-  
355 uals showing higher than 50 % of wave contribution, but in older system of Spain, it is not  
356 the case. A possible speculation is that, in the process of speciation, hybrids are more likely to  
357 mate with crab (or to say bigger individuals), then hybrid population in PCA become closer and  
358 closer to crab population, ancestry proportion of wave decrease and proportion of crab increase.  
359 In addition of assortative mating, the possible reason of one direction of introgression could also  
360 be asymmetrical selection, or confusion of adaptation to the lower shore within crab group or a  
361 mix of these.

362 **4.2 Summary statistics of differentiation**

363 Here I expect to find special population genetic differentiation and nucleotide density pattern  
364 in chromosomes, detect new inversions candidates and infer parallel adaptive evolution history.

365 **4.2.1 SFS**

366 The 1d SFS seems to be far from the neutral expectation. It could not be caused by introgression  
367 since it is similar in all groups. The SFS might be distorted by SNPs in polymorphic inversions.  
368 I only used sites detected in all individuals to guarantee accuracy of gene frequency, but this  
369 might also distort the SFS. For instance, if repeat regions are more likely to fall into this class,  
370 the SFS will be inaccurate. From 2d SFS, the only significant information is that crab and  
371 hybrid populations are closer which is consistent with PCA result.

372 **4.2.2 Population genetic differentiation**

373 Most of loci are not fixed and keep polymorphic. Rare fixed loci suggest balancing selection  
374 which means heterozygotes have a higher fitness than homozygotes. The same pattern for SNPs  
375 was found by Westram (Westram et al. 2018) and some possible explanations were given: in-  
376 direct divergent selection on SNPs linked to selected variants, selection on polygenic traits or  
377 a combination of divergent selection and balancing selection that make loci maintain polymor-  
378 phism in one or two habitat ends.

379

380 High FST means high divergence. Regions with lots of high FST values may well contain genes  
381 under divergent selection between crab and wave ecotypes although along with some noises. It  
382 is hard to find these regions just from FST scatter plot due to some regions' more SNPs. Usu-  
383 ally, regions with many SNPs tend to own high FST SNPs. Therefore, I colored top 0.01% of  
384 SNPs as Figure5A and calculated mean FST per map position showed as Figure5B. There are  
385 some of the same regions detected before by Morales (Morales et al. 2019) including inversion  
386 candidates. There are also some regions on LG2, 8, 11 demonstrating high population genetic  
387 differentiation but not regarded as inversion candidates. However, further exploration should  
388 be executed to verify whether they are new inversions or just genetic speciation islands which  
389 harbor loci underlying reproductive isolation.

390

391 Generally, I expect to find high FST, high PBS, low pi and high LD region as candidate blocks  
392 for reproductive isolation. Chromosome inversion candidates are supposed to carry many loci  
393 under divergent selection, so I expect to detect the same summary statistics pattern in these  
394 candidates. However, I found that many candidates show higher pi, even lower FST and PBS.  
395 Only LG1.2, 2.1, 4.1, 6.1, 6.2, 9.1 meet all conditions. Higher pi is possible when inversions are  
396 polymorphic. The elevated diversity could be between arrangements or within an arrangement  
397 at the nucleotide level due to balancing selection. For example, LG10.1, 12.1, 12.2, 14.1, 14.2  
398 show lower mean FST, lower mean PBS but higher mean pi and higher mean LD. Balancing  
399 selection may also result in lower FST and lower PBS, but these regions still maintain high  
400 LD. Some inversion regions have similar levels of diversity to whole chromosome suggesting  
401 non-polymorphic and not accumulated differentiation. Inversion regions have less diversity sug-  
402 gesting these could have swept to high frequency recently.

403

404 It is worth noting that chromosome inversions can't fully explain parallelism. Morales also found  
405 shared outlier loci distributing across the genome (Morales et al. 2019). This could be resulted  
406 by polygenic selection of multiple loci with small effect underlying parallel adaptive divergence.

407 **4.2.3 Limitation of FST**

408 In genomic regions surrounding barrier loci, the barrier effect initially allow a build-up of genetic  
409 differentiation in the form of allele frequency variation between populations, typically using a  
410 relative measure such as FST (Ravinet et al. 2017). However, Using FST as indicator of diver-  
411 gent selection has some problems.

412

413 The pattern of FST is affected by multiple factors that change throughout the genome, including  
414 mutation, genetic drift, selection, demographic history, recombination, gene flow, gene density,  
415 and genome structure, some of which are expected to change at different stages of speciation.  
416 Polygenic makes it possible that FST of many underlying loci remain low when snails achieving  
417 trait divergence(Westram et al. 2014). Background and positive selection may produce similar  
418 patterns in genome scan using FST, reducing intraspecific diversity and increasing interspecific  
419 gene composition(Cruickshank & Hahn 2014). Furthermore, some non-inversion candidates  
420 regions with high FST may not be other currently undetected chromosomal inversions under  
421 positive divergence selection, but related to other mechanisms suppressing recombination such  
422 as centromere region combined with background selection or divergence hitchhiking (Ravinet  
423 et al. 2017).

424 **4.3 Linkage disequilibrium**

425 From LD analysis, I could find that in 17 chromosomes, crab population always has the highest  
426 LD value. LD is likely to be high in crab for many reasons. It might because crab has a small  
427 effective population size, or crab suffered bottleneck accident in history, or crab is under a rela-  
428 tively strong selection, or crab individuals have a stronger tend of assortative mating. Genetic  
429 structure may also increase LD. Polymorphic inversions are able to suppress recombination.  
430 The clines generate LD whether they are due to introgression or due to a selective gradient.  
431 Generally, hybrid has medium LD value and wave has the lowest LD value, but the difference  
432 is little. Selection of wave environment is not as strong as crab, so it is reasonable that LD  
433 value of wave ecotype population is the lowest, and hybrid is medium. The little LD difference  
434 between hybrid and wave group indicates that the directional introgression destroys high LD  
435 genetic structure in crab group.

436

437 LD heat map should reveal inversions and a general pattern of higher LD among neighboring  
438 loci, especially among SNPs in the same contig. However, in LD heat map of 17 chromosomes in  
439 *Littorina saxatilis*, there is usually no significant higher LD value block. Chromosome inversion  
440 candidates LG4.1, 7.2, 9.1, 11.1, 14.1, 14.2 are relatively clear to find in whole chromosome LD  
441 heat map compared with other inversion candidates. Therefore, it is more difficult to infer new  
442 inversion candidates just from LD heatmap.

443

444 For LD decay analysis, crab usually has the slowest LD decay rate in 17 chromosomes which  
445 means selection effect on crab ecotype population is stronger than other population. Selection  
446 will decrease population genetic diversity and strengthen association between pairwise loci (LD).  
447 Therefore, usually population under stronger selection has slower LD decay rate. *Littorina* is  
448 a high fertility and short generation interval species, therefore LD decay rate in this species  
449 is generally fast. Hybrid ecotype population LD decay rate is usually the fastest which might  
450 be resulted by higher population genetic diversity gained from crab and wave ecotype population.

451

452 Points disperse along the distance may be resulted by polymorphic inversions on chromosomes  
453 which makes high LD between SNPs far apart. In some chromosomes, points disperse more early  
454 than others which might be because short inversions on chromosome results in LD between loci  
455 far away. The most obvious effect of inversions on LD decay is typically going to be on larger  
456 scales especially between breakpoints far apart (such as 10Mb), but due to limitation of contig  
457 size, I could not find this situation.

458

459 Ravinet (Ravinet et al. 2016) found that LD clusters showed strong signals of ecotype specificity,  
460 loci were largely homozygous in crab while heterozygous in wave, suggesting these loci have  
461 undergone selective sweeps (the reduction or elimination of variation at sites that are physically  
462 linked to a site under directional selection) in crab population. This is another proof that crab  
463 group suffered stronger selection fueling divergence.

#### 464 4.4 Other possible mechanisms of parallel adaptation

465 Morales (Morales et al. 2019) supposed that standing variation within chromosomal inversions  
466 can be maintained as balanced polymorphism and fuel rapid parallel phenotypic divergence to  
467 heterogeneous environments through gene flow without high sharing of genomic outliers. Bal-  
468 ancing and divergent selection between habitats could maintain inversions for long periods of  
469 time, resulting in the high diversity and divergence for some inversions noted above (Faria  
470 et al. 2019). Furthermore, balancing selection is often documented for inversion polymorphisms  
471 (Wellenreuther & Bernatchez 2018). Inversions can extend the impact of barrier loci (reproduc-  
472 tive isolation loci) to linked loci and promote additional barrier loci accumulating in inverted  
473 region with gene flow between populations which facilitates the efficient spread of adaptive  
474 standing variance in inversions (Morjan & Rieseberg 2004). After the establishment of local  
475 LD, gene drift can facilitate disruptive selection (Dieckmann & Doebeli 1999) which suggests  
476 random process may play a role in speciation.

477

478 There are some study revealing that a large scale of SNPs for some key adaptive phenotypes  
479 or environmental axes could be explained by SNPs in chromosome inversions (Morales et al.  
480 2019). Westram also found that 75% of non-neutral SNPs cluster together in three LGs (6,  
481 14, 17) which suggest that large regions of low recombination consistent with the presence of  
482 chromosomal rearrangements.

#### 483 4.5 Future direction

484 From a speciation perspective, the main goal is to infer the number, distribution and intensity  
485 of gene flow barriers, and their impact on other genomic regions, then find genetic speciation is-  
486 land which harbor loci underlying reproductive isolation. However, due to many other processes'  
487 existence, it is challengeable to detect this signal from genome scan. Chromosome inversions  
488 are regarded as candidates for gene flow barriers. Then the research on ages and spatial distri-  
489 butions of inversions will be important to understand local adaptation and evolution history of  
490 reproductive isolation in *Littorina saxatilis*. It is premature to use diversity and divergence data

491 of *Littorina saxatilis* to calculate exact inversion ages, but Morals (Morales et al. 2019) inferred  
492 from Sweden snails data that the origins of inversions are earlier than postglacial colonization  
493 of the Swedish coast.

494

495 After finding inversions region under divergence selection, the next step is to identify why they  
496 are so differentiated. These regions are likely to contain reproductive isolation loci which need  
497 more exploration, such as finding break point of inversions, functional annotation of outlier loci,  
498 selection scans, phenotype-habitat-genotype associations, experimental analyses with sequenc-  
499 ing. It is worth mentioning that size is the feature under strong divergence selection and also  
500 could be main reason of assortative mating between ecotypes (Smadja & Butlin 2011), therefore  
501 size is likely to facilitate the formation of barriers between ecotypes. Then finding the associa-  
502 tion between size and reproductive isolation which is deduced by Johannesson (Johannesson  
503 2016) is another promising research direction.

504 **5 Conclusion**

505 In general, I found two distinct groups in the Spanish system of *Littorina saxatilis* instead of a  
506 cline displayed in the more recent system Sweden. Compared with the wave ecotype, hybrids are  
507 much closer to the crab ecotype. Besides, the directional introgression from wave to crab destroys  
508 high LD in the crab ecotype. Differentiation analysis hints at balancing selection on this species  
509 and LD analysis hints at stronger selection and more polymorphic inversions of crab ecotype.  
510 Also, I found some new chromosome inversions candidates for future verification. This research  
511 tried to infer the process and mechanisms of ecological divergence and speciation after the phase  
512 of the Sweden system. However, to explore the underlying mechanisms of parallel ecological  
513 divergence with gene flow and the evolvement of reproductive isolation with interbreeding, it  
514 still needs more effort in future.

515 **Code availability**

516 All of the scripts, software parameter settings, and dependancies packages are included in the  
517 following website: <https://github.com/rz520/LittorinaPipeline>

## 518 Acknowledgement

519 Here I would like to appreciate supervisor Matteo Fumagalli for his continuous help on software  
520 usage, parameter setting, and results interpretation. Co-supervisor Francesca Raffini and Roger  
521 K. Butlin gave me lots of guidance on *Littorina saxatilis* background knowledge and recent  
522 research progress. All of my supervisors helped me to polish my thesis by giving constructive  
523 suggestion. I am also grateful to other master students in Matteo's Lab and my friends for their  
524 help, encouragement and company.

525 **References**

- 526 Beaumont, M. A. (2010), ‘Approximate bayesian computation in evolution and ecology’, *Annual  
527 Review of Ecology, Evolution, and Systematics* **41**(1), 379–406.  
528 **URL:** <https://doi.org/10.1146/annurev-ecolsys-102209-144621>
- 529 Butlin, R. K., Galindo, J. & Grahame, J. W. (2008), ‘Sympatric, parapatric or allopatric: the  
530 most important way to classify speciation?’, *Philosophical Transactions of the Royal Society  
531 B: Biological Sciences* **363**(1506), 2997–3007.  
532 **URL:** <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.2008.0076>
- 533 Cruickshank, T. E. & Hahn, M. W. (2014), ‘Reanalysis suggests that genomic islands of specia-  
534 tion are due to reduced diversity, not reduced gene flow’, *Molecular Ecology* **23**(13), 3133–3157.  
535 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.12796>
- 536 Dieckmann, U. & Doebeli, M. (1999), On the origin of species by sympatric speciation, Iiasa  
537 interim report, IIASA, Laxenburg, Austria.  
538 **URL:** <http://pure.iiasa.ac.at/id/eprint/5926/>
- 539 Faria, R., Chaube, P., Morales, H. E., Larsson, T., Lemmon, A. R., Lemmon, E. M., Rafajlović,  
540 M., Panova, M., Ravinet, M., Johannesson, K., Westram, A. M. & Butlin, R. K. (2019),  
541 ‘Multiple chromosomal rearrangements in a hybrid zone between littorina saxatilis ecotypes’,  
542 *Molecular Ecology* **28**(6), 1375–1393.  
543 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.14972>
- 544 Fitzpatrick, B. M., Fordyce, J. A. & Gavrilets, S. (2008), ‘What, if anything, is sympatric  
545 speciation?’, *Journal of Evolutionary Biology* **21**(6), 1452–1459.  
546 **URL:** <https://doi.org/10.1111/j.1420-9101.2008.01611.x>
- 547 Fox, E. A., Wright, A. E., Fumagalli, M. & Vieira, F. G. (2019), ‘ngsLD: evaluating linkage  
548 disequilibrium using genotype likelihoods’, *Bioinformatics* **35**(19), 3855–3856.  
549 **URL:** <https://doi.org/10.1093/bioinformatics/btz200>
- 550 Galindo, J., Martínez-Fernández, M., Rodríguez-Ramilo, S. T. & Rolán-Alvarez, E. (2013),  
551 ‘The role of local ecology during hybridization at the initial stages of ecological speciation in  
552 a marine snail’, *Journal of Evolutionary Biology* **26**(7), 1472–1487.  
553 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/jeb.12152>
- 554 Galindo, J., Morán, P. & Rolán-Alvarez, E. (2009), ‘Comparing geographical genetic differentia-  
555 tion between candidate and noncandidate loci for adaptation strengthens support for parallel  
556 ecological divergence in the marine snail littorina saxatilis’, *Molecular Ecology* **18**(5), 919–930.  
557 **URL:** <https://doi.org/10.1111/j.1365-294X.2008.04076.x>
- 558 Gilly, A., Southam, L., Suveges, D., Kuchenbaecker, K., Moore, R., Melloni, G. E. M., Hatziko-  
559 toulas, K., Farmaki, A.-E., Ritchie, G., Schwartzentuber, J., Danecek, P., Kilian, B., Pollard,

- 560 M. O., Ge, X., Tsafantakis, E., Dedoussis, G. & Zeggini, E. (2018), 'Very low-depth whole-  
561 genome sequencing in complex trait association studies', *Bioinformatics* **35**(15), 2555–2561.  
562 **URL:** <https://doi.org/10.1093/bioinformatics/bty1032>
- 563 Grahame, J. W., Wilding, C. S. & Butlin, R. K. (2006), 'Adaptation to a steep environ-  
564 mental gradient and an associated barrier to gene exchange in littorina saxatilis', *Evolution*  
565 **60**(2), 268–278.  
566 **URL:** <https://doi.org/10.1111/j.0014-3820.2006.tb01105.x>
- 567 Hollander, J., Lindegarth, M. & Johannesson, K. (2005), 'Local adaptation but not geographical  
568 separation promotes assortative mating in a snail', *Animal Behaviour* **70**(5), 1209–1219.  
569 **URL:** <https://www.sciencedirect.com/science/article/pii/S0003347205002733>
- 570 Johannesson, K. (2016), 'What can be learnt from a snail?', *Evolutionary Applications* **9**(1), 153–  
571 165.  
572 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/eva.12277>
- 573 Johannesson, K., Johannesson, B. & Rolán-Alvarez, E. (1993), 'Morphological differentiation  
574 and genetic cohesiveness over a microenvironmental gradient in the marine snail littorina  
575 saxatilis', *Evolution* **47**(6), 1770–1787.  
576 **URL:** <https://doi.org/10.1111/j.1558-5646.1993.tb01268.x>
- 577 Johannesson, K., Panova, M., Kemppainen, P., André, C., Rolán-Alvarez, E. & Butlin, R. K.  
578 (2010), 'Repeated evolution of reproductive isolation in a marine snail: unveiling mecha-  
579 nisms of speciation', *Philosophical Transactions of the Royal Society B: Biological Sciences*  
580 **365**(1547), 1735–1747.  
581 **URL:** <https://royalsocietypublishing.org/doi/abs/10.1098/rstb.2009.0256>
- 582 Johannesson, K., Rolán-Alvarez, E. & Ekendahl, A. (1995), 'Incipient reproductive isolation  
583 between two sympatric morphs of the intertidal snail littorina saxatilis', *Evolution* **49**(6), 1180–  
584 1190.  
585 **URL:** <https://doi.org/10.1111/j.1558-5646.1995.tb04445.x>
- 586 Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. (2014), 'ANGSD: Analysis of next generation  
587 sequencing data', *BMC Bioinformatics* **15**(1), 356.  
588 **URL:** <http://www.biomedcentral.com/1471-2105/15/356/abstract>
- 589 Meisner, J. & Albrechtsen, A. (2018), 'Inferring population structure and admixture proportions  
590 in low-depth ngs data', *Genetics* **210**(2), 719–731.  
591 **URL:** <https://www.genetics.org/content/210/2/719>
- 592 Morales, H. E., Faria, R., Johannesson, K., Larsson, T., Panova, M., Westram, A. M. & Butlin,  
593 R. K. (2019), 'Genomic architecture of parallel ecological divergence: Beyond a single envi-  
594 ronmental contrast', *Science Advances* **5**(12).  
595 **URL:** <https://advances.sciencemag.org/content/5/12/eaav9963>

- 596 Morjan, C. L. & Rieseberg, L. H. (2004), 'How species evolve collectively: implications of gene  
597 flow and selection for the spread of advantageous alleles', *Molecular Ecology* **13**(6), 1341–1356.  
598 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1365-294X.2004.02164.x>
- 599 Panova, M., Aronsson, H., Cameron, R., Dahl, P., Godhe, A., Lind, U., Ortega-Martinez,  
600 O., Pereyra, R., Tesson, S., Wrangé, A.-L., Blomberg, A. & Johannesson, K. (2016), *DNA  
601 Extraction Protocols for Whole-Genome Sequencing in Marine Organisms*, Vol. 1452, pp. 13–  
602 44.  
603 **URL:** <https://pubmed.ncbi.nlm.nih.gov/27460368/>
- 604 Panova, M., Blakeslee, A. M. H., Miller, A. W., Mäkinen, T., Ruiz, G. M., Johannesson, K.  
605 & André, C. (2011), 'Glacial history of the north atlantic marine snail, littorina saxatilis,  
606 inferred from distribution of mitochondrial dna lineages', *PLOS ONE* **6**(3), 1–14.  
607 **URL:** <https://doi.org/10.1371/journal.pone.0017511>
- 608 Panova, M., Hollander, J. & Johannesson, K. (2006), 'Site-specific genetic divergence in parallel  
609 hybrid zones suggests nonallopatric evolution of reproductive barriers', *Molecular Ecology*  
610 **15**(13), 4021–4031.  
611 **URL:** <https://doi.org/10.1111/j.1365-294X.2006.03067.x>
- 612 Rafajlovic, M., Eriksson, A., Rimark, A., Hintz-Saltin, S., Charrier, G., Panova, M., Andre, C.,  
613 Johannesson, K. & Mehlig, B. (2013), 'The Effect of Multiple Paternity on Genetic Diversity  
614 of Small Populations during and after Colonisation', *PLOS ONE* **8**(10).  
615 **URL:** <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0075587>
- 616 Ravinet, M., Faria, R., Butlin, R. K., Galindo, J., Bierne, N., Rafajlović, M., Noor, M. A. F.,  
617 Mehlig, B. & Westram, A. M. (2017), 'Interpreting the genomic landscape of speciation: a  
618 road map for finding barriers to gene flow', *Journal of Evolutionary Biology* **30**(8), 1450–1477.  
619 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/jeb.13047>
- 620 Ravinet, M., Westram, A., Johannesson, K., Butlin, R., André, C. & Panova, M. (2016), 'Shared  
621 and nonshared genomic divergence in parallel ecotypes of littorina saxatilis at a local scale',  
622 *Molecular Ecology* **25**(1), 287–305.  
623 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.13332>
- 624 Reid, D. G. (1996), *Systematics and evolution of Littorina*, Ray Society.  
625 **URL:** <https://www.nhbs.com/systematics-and-evolution-of-littorina-book>
- 626 Rolan-Alvarez, E., Johannesson, K. & Erlandsson, J. (1997), 'The maintenance of a cline in  
627 the marine snail Littorina saxatilis: The role of home site advantage and hybrid fitness',  
628 *EVOLUTION* **51**(6), 1838–1847.  
629 **URL:** <https://doi.org/10.1111/j.1558-5646.1997.tb05107.x>
- 630 Schlüter, D. (2009), 'Evidence for ecological speciation and its alternative', *Science*  
631 **323**(5915), 737–741.  
632 **URL:** <https://science.sciencemag.org/content/323/5915/737>

- 633 Skotte, L., Korneliussen, T. S. & Albrechtsen, A. (2013), 'Estimating Individual Admixture  
634 Proportions from Next Generation Sequencing Data', *Genetics* **195**(3), 693–702.  
635 **URL:** <https://doi.org/10.1534/genetics.113.154138>
- 636 Smadja, C. M. & Butlin, R. K. (2011), 'A framework for comparing processes of speciation in  
637 the presence of gene flow', *MOLECULAR ECOLOGY* **20**(24), 5123–5140.  
638 **URL:** <https://doi.org/10.1111/j.1365-294X.2011.05350.x>
- 639 Therkildsen, N. O., Wilder, A. P., Conover, D. O., Munch, S. B., Baumann, H. & Palumbi,  
640 S. R. (2019), 'Contrasting genomic shifts underlie parallel phenotypic evolution in response  
641 to fishing', *Science* **365**(6452), 487–490.  
642 **URL:** <https://science.sciencemag.org/content/365/6452/487>
- 643 Wellenreuther, M. & Bernatchez, L. (2018), 'Eco-evolutionary genomics of chromosomal inver-  
644 sions', *Trends in Ecology & Evolution* **33**(6), 427–440.  
645 **URL:** <https://www.sciencedirect.com/science/article/pii/S0169534718300788>
- 646 Westram, A. M., Galindo, J., Alm Rosenblad, M., Grahame, J. W., Panova, M. & Butlin,  
647 R. K. (2014), 'Do the same genes underlie parallel phenotypic divergence in different littorina  
648 saxatilis populations?', *Molecular Ecology* **23**(18), 4603–4616.  
649 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.12883>
- 650 Westram, A. M., Rafajlović, M., Chaube, P., Faria, R., Larsson, T., Panova, M., Ravinet, M.,  
651 Blomberg, A., Mehlig, B., Johannesson, K. & Butlin, R. (2018), 'Clines on the seashore: The  
652 genomic architecture underlying rapid divergence in the face of gene flow', *Evolution Letters*  
653 **2**(4), 297–309.  
654 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1002/evl3.74>
- 655 Wilder, A. P., Palumbi, S. R., Conover, D. O. & Therkildsen, N. O. (2020), 'Footprints of local  
656 adaptation span hundreds of linked genes in the atlantic silverside genome', *Evolution Letters*  
657 **4**(5), 430–443.  
658 **URL:** <https://onlinelibrary.wiley.com/doi/abs/10.1002/evl3.189>
- 659 Yeaman, S. (2013), 'Genomic rearrangements and the evolution of clusters of locally adaptive  
660 loci', *Proceedings of the National Academy of Sciences* **110**(19), E1743–E1751.  
661 **URL:** <https://www.pnas.org/content/110/19/E1743>
- 662 Zhang, C., Dong, S.-S., Xu, J.-Y., He, W.-M. & Yang, T.-L. (2018), 'PopLDdecay: a fast  
663 and effective tool for linkage disequilibrium decay analysis based on variant call format files',  
664 *Bioinformatics* **35**(10), 1786–1788.  
665 **URL:** <https://doi.org/10.1093/bioinformatics/bty875>

666 Appendix

Linkage group (LG)	LD cluster	Cluster size (cM)	Start (cM)	End (cM)	Mean FST	Mean PBS	Mean Pi	Mean LD (r2)	candidate
LG1	LGC1.1	2.1	0	2.1	0.0019398	0.012677	0.00648093	0.1939411	1
LG1	LGC1.2	5.42	75.53	80.95	0.00219029	0.02014028	0.00536197	0.2134792	1
LG1	LG1	80.95	0	80.95	0.00212763	0.01073747	0.00594011	0.1704242	0
LG2	LGC2.1	13.87	0.34	14.21	0.00614239	0.01879701	0.00576577	0.1999839	1
LG2	LG2	88.76	0	88.76	0.00274161	0.01610663	0.00592564	0.1646505	0
LG3	LG3	68.02	0	68.02	0.00288605	0.01137463	0.00603668	0.1761798	0
LG4	LGC4.1	0.48	1.03	1.51	0.00463436	0.03770791	0.00587387	0.1847306	1
LG4	LG4	56.52	0	56.52	0.00243	0.0140493	0.00587649	0.1711099	0
LG5	LG5	54.07	0	54.07	0.00280035	0.01904912	0.00591619	0.1699921	0
LG6	LGC6.1	29.3	0	29.3	0.00297457	0.03524753	0.00567617	0.1820881	1
LG6	LGC6.2	20.57	8.73	29.3	0.00269045	0.02541398	0.00579544	0.1825618	1
LG6	LG6	60.25	0	60.25	0.00227857	0.01616375	0.00634677	0.1691389	0
LG7	LGC7.1	1.73	36.01	37.74	0.00275902	0.01656839	0.00683345	0.2117835	1
LG7	LGC7.2	9.29	42.08	51.37	0.00364395	0.02627485	0.00641327	0.1799024	1
LG7	LG7	51.37	0	51.37	0.00235047	0.0124258	0.00629948	0.1714745	0
LG8	LG8	54.38	0	54.38	0.0025284	0.01993345	0.00671042	0.186891	0
LG9	LGC9.1	23.18	18.64	41.82	0.00275828	0.0234772	0.00566658	0.1917811	1
LG9	LG9	56.67	0	56.67	0.00256584	0.02090745	0.00763238	0.1833917	0
LG10	LGC10.1	2.54	0.58	3.12	0.00251415	0.01067118	0.00638226	0.1959969	1
LG10	LG10	45.53	0	45.53	0.00280773	0.02139548	0.00636286	0.188335	0
LG11	LGC11.1	0.59	52.32	52.91	0.00259493	0.0145582	0.00686731	0.2100456	1
LG11	LG11	58.39	0	58.39	0.0024994	0.02136298	0.00675471	0.1950033	0
LG12	LGC12.1	26.31	3.32	29.63	0.00155347	0.01451515	0.00667943	0.1775293	1
LG12	LGC12.2	11.52	48.71	60.24	0.00299152	0.01851379	0.00654778	0.1870268	1
LG12	LG12	60.24	0	60.24	0.00314249	0.02937913	0.00609529	0.1726986	0
LG13	LG13	51.3	0	51.3	0.0018803	0.01141404	0.00618406	0.1804347	0
LG14	LGC14.1	11.32	0.39	11.71	0.00325176	0.00927283	0.0060404	0.1966089	1
LG14	LGC14.2	2.9	8.81	11.71	0.00286758	0.00959897	0.00606569	0.1960541	1
LG14	LGC14.3	23.23	11.71	34.94	0.0034732	0.01206736	0.00528926	0.1782203	1
LG14	LG14	48.71	0	48.71	0.00332075	0.00996925	0.00583716	0.1881616	0
LG15	LG15	56.49	0	56.49	0.00201495	0.00941298	0.00633721	0.1755119	0
LG16	LG16	57.44	0	57.44	0.00212438	0.01346309	0.00695419	0.1911745	0
LG17	LGC17.1	15.33	46.99	62.32	0.00561686	0.0379331	0.00655994	0.1748176	1
LG17	LG17	62.32	0	62.32	0.00219622	0.01432516	0.00598865	0.1672243	0

Table 1: Mean summary statistics in inversion candidates and whole chromosome

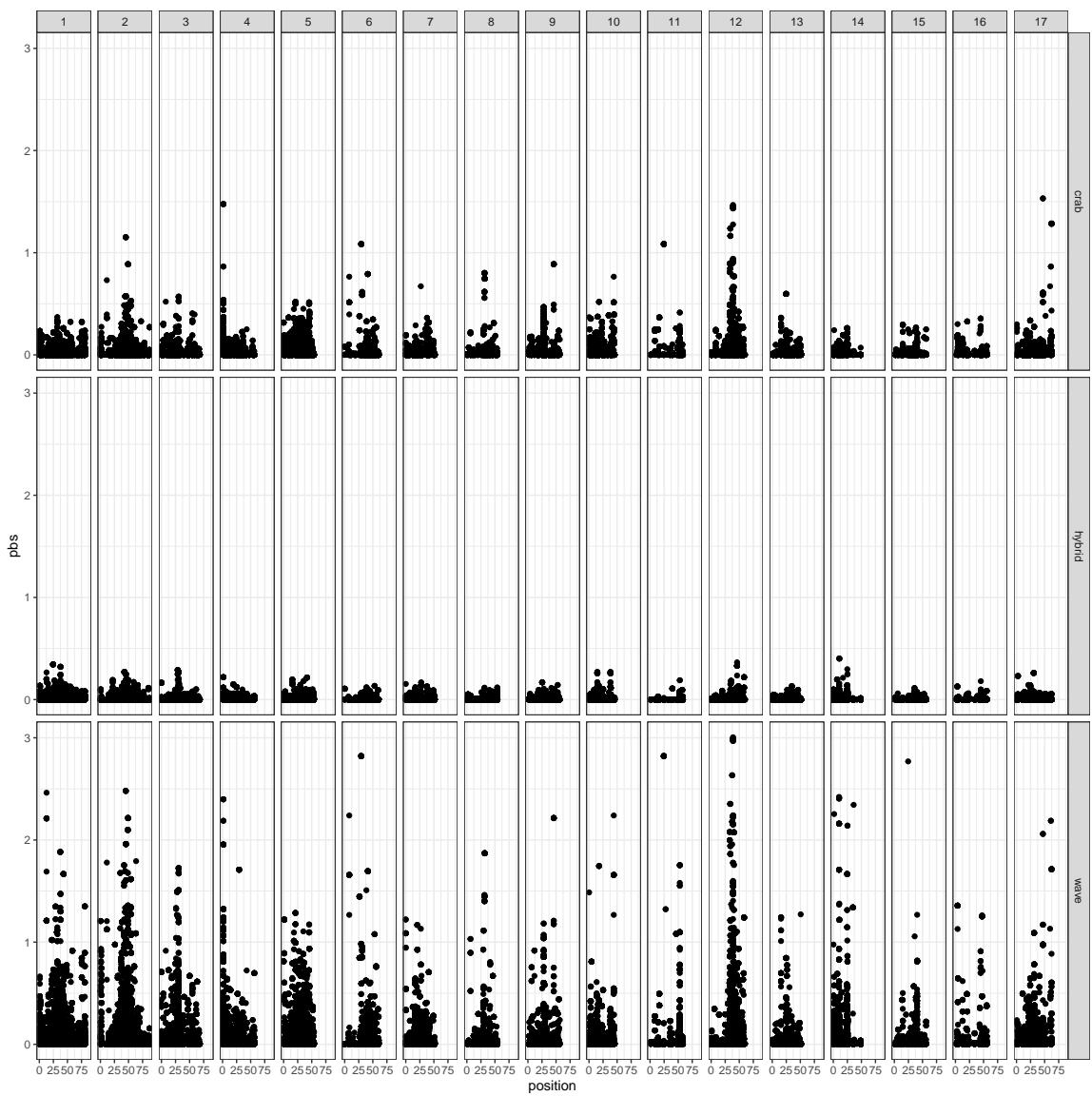


Figure 1: PBS value scatter plot of every chromosome and population in windows

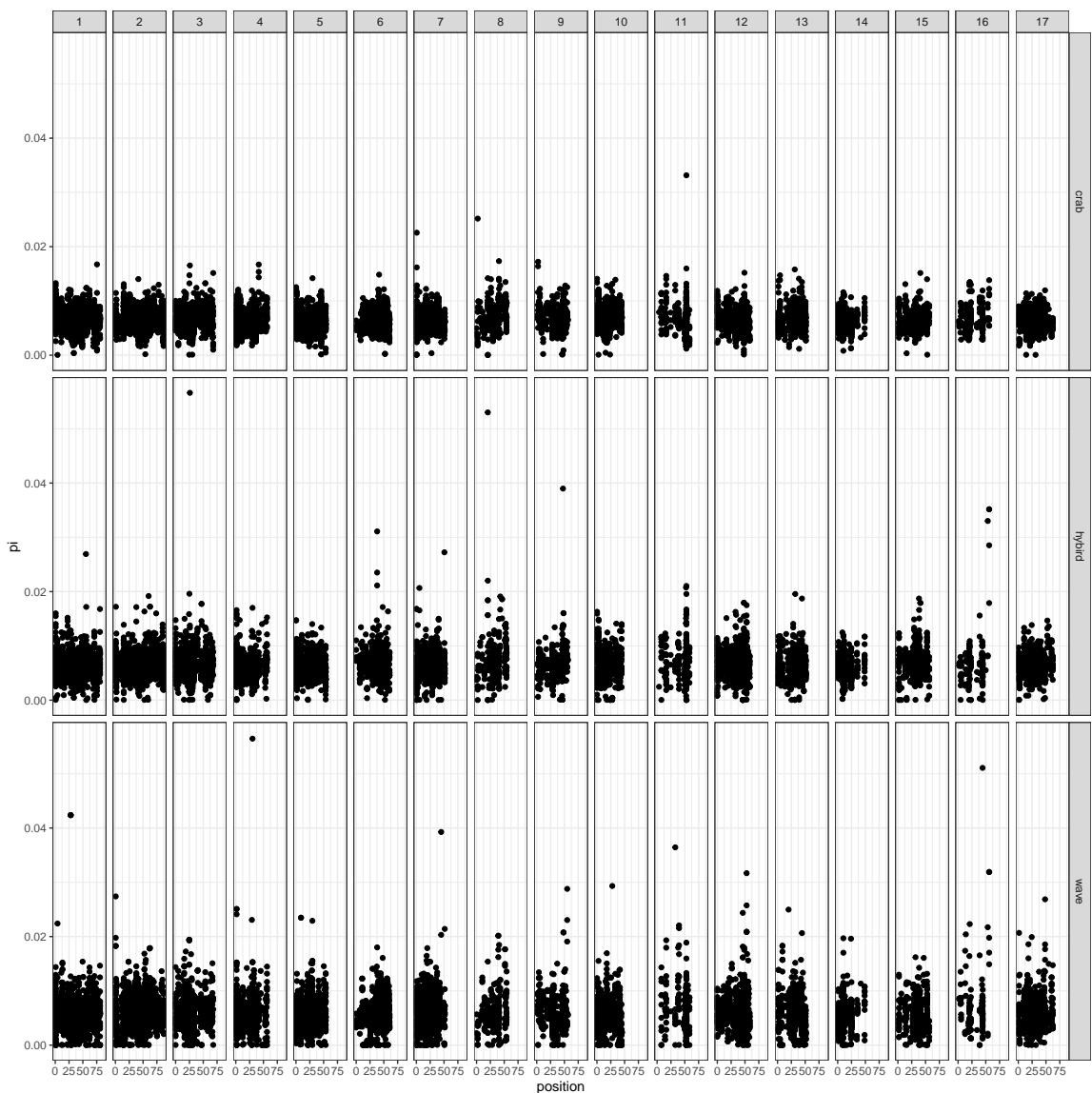


Figure 2: Pi value scatter plot of every chromosome and population in windows

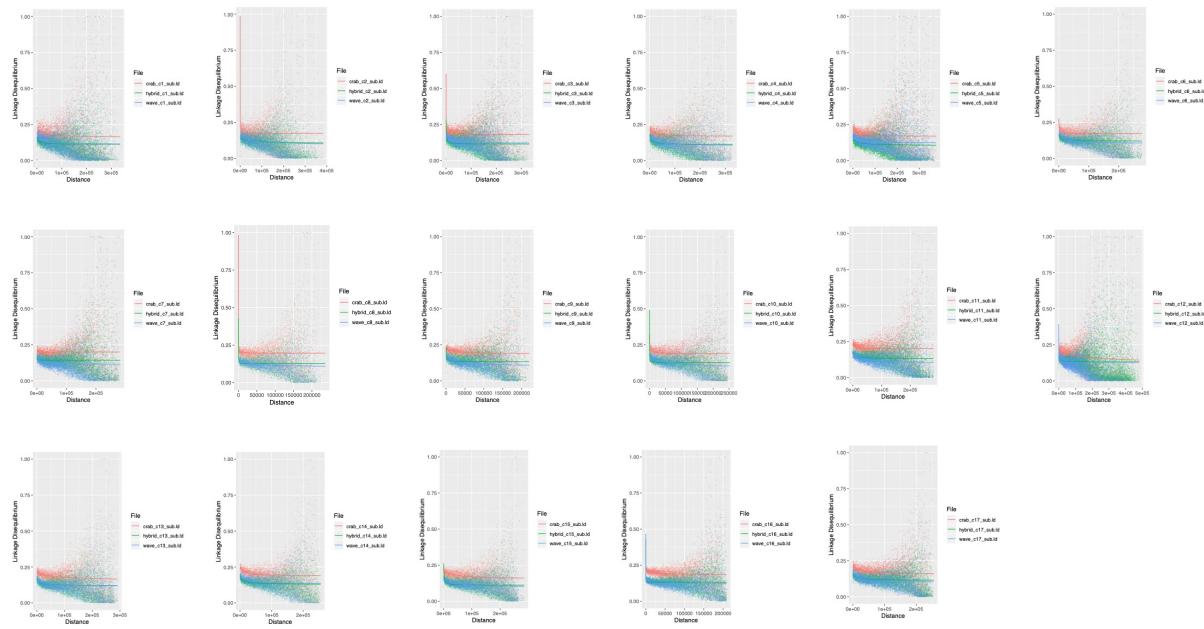


Figure 3: LD decay analysis of every chromosome

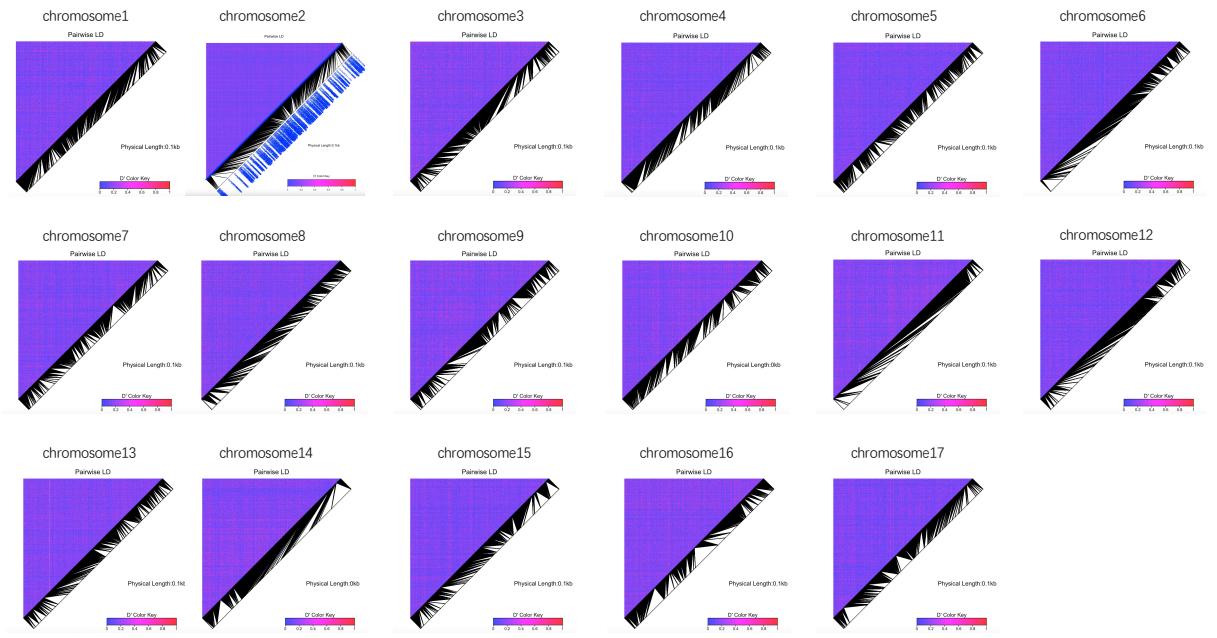


Figure 4: LD heat map of every chromosome